

Audio classification using machine learning

Data: <https://urbansounddataset.weebly.com/download-urbansound8k.html> (<https://urbansounddataset.weebly.com/download-urbansound8k.html>)

```
In [1]: 1 # !pip install librosa
```

```
In [2]: 1 import numpy as np
        2
```

```
In [3]: 1 import tensorflow as tf
```

```
In [4]: 1 try:
        2     import librosa
        3     print("Librosa is available.")
        4 except ImportError:
        5     print("Librosa is not installed.")
        6
```

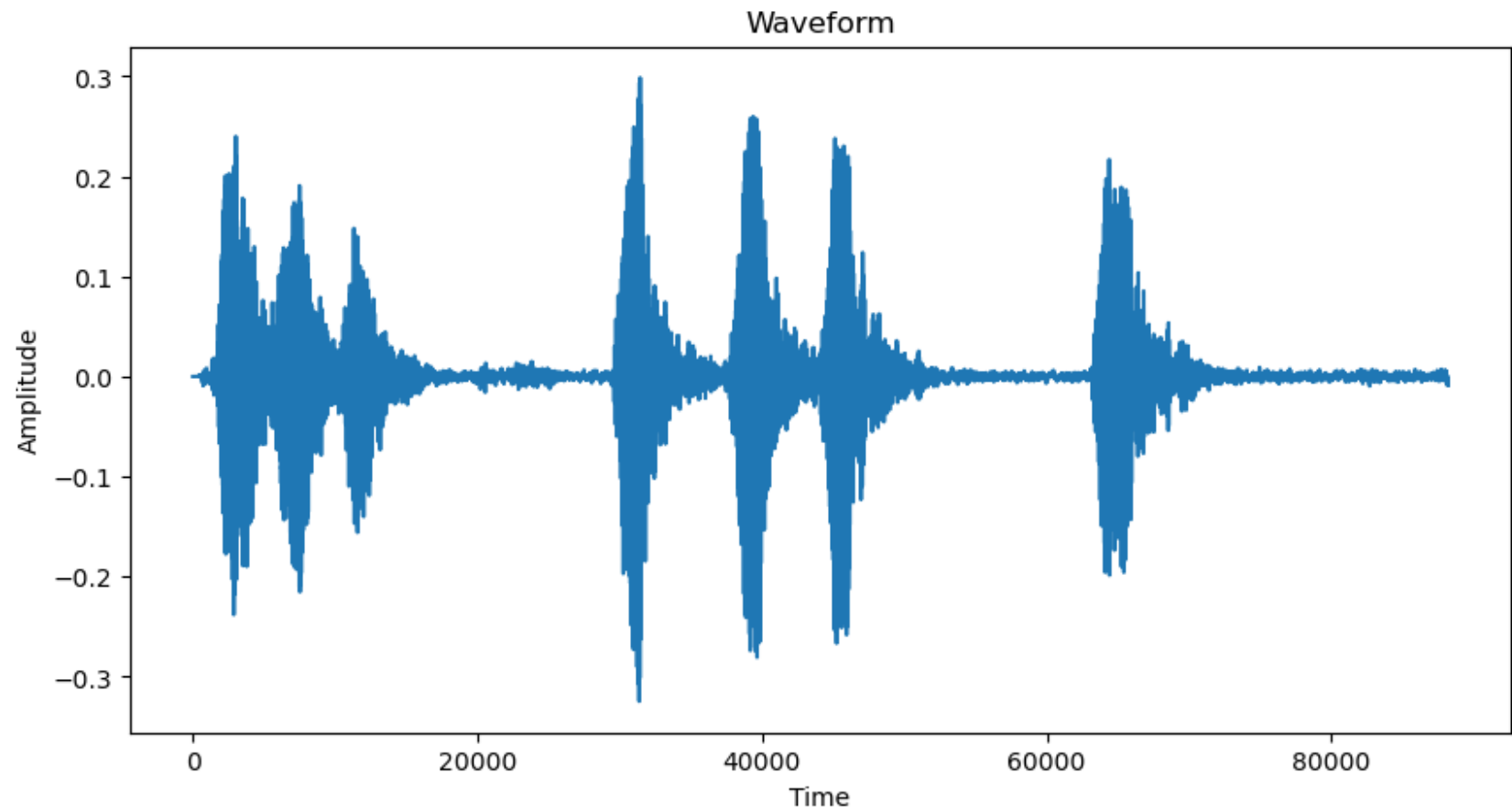
Librosa is available.

```
In [5]: 1 import matplotlib.pyplot as plt
        2 %matplotlib inline
```

```
In [6]: 1 import IPython.display as ipd
        2 import librosa
        3 import librosa.display
```

```
In [7]: 1 filename = 'dog bark.wav'
```

```
In [8]: 1 data, sample_rate = librosa.load(filename)
2 plt.figure(figsize=(10, 5))
3 plt.plot(data)
4 plt.title('Waveform')
5 plt.xlabel('Time')
6 plt.ylabel('Amplitude')
7 plt.show()
8
9 ipd.Audio(filename)
```



Out[8]:

0:04 / 0:04

```
In [9]: 1 sample_rate
```

```
Out[9]: 22050
```

```
In [10]: 1 from scipy.io import wavfile as wav
2 wave_sample_rate, wave_audio = wav.read(filename)
```

```
In [11]: 1 wave_sample_rate
```

```
Out[11]: 44100
```

```
In [12]: 1 wave_audio
```

```
Out[12]: array([[ 0,  0],
               [ 0,  0],
               [ 0,  0],
               ...,
               [-399, -115],
               [-388, -111],
               [-386, -105]], dtype=int16)
```

```
In [13]: 1 data
```

```
Out[13]: array([ 3.4924597e-10,  3.4924597e-10,  4.6566129e-10, ...,
                -7.9498515e-03, -7.7366987e-03, -8.0531817e-03], dtype=float32)
```

```
In [14]: 1 import pandas as pd
2
3 metadata = pd.read_csv('UrbanSound8K/metadata/UrbanSound8K.csv')
```

In [15]: 1 metadata.head(10)

Out[15]:

	slice_file_name	fsID	start	end	salience	fold	classID	class
0	100032-3-0-0.wav	100032	0.000000	0.317551	1	5	3	dog_bark
1	100263-2-0-117.wav	100263	58.500000	62.500000	1	5	2	children_playing
2	100263-2-0-121.wav	100263	60.500000	64.500000	1	5	2	children_playing
3	100263-2-0-126.wav	100263	63.000000	67.000000	1	5	2	children_playing
4	100263-2-0-137.wav	100263	68.500000	72.500000	1	5	2	children_playing
5	100263-2-0-143.wav	100263	71.500000	75.500000	1	5	2	children_playing
6	100263-2-0-161.wav	100263	80.500000	84.500000	1	5	2	children_playing
7	100263-2-0-3.wav	100263	1.500000	5.500000	1	5	2	children_playing
8	100263-2-0-36.wav	100263	18.000000	22.000000	1	5	2	children_playing
9	100648-1-0-0.wav	100648	4.823402	5.471927	2	10	1	car_horn

In [16]: 1 # check wether the dataset is imbalanced
2 metadata['class'].value_counts()

Out[16]: dog_bark 1000
children_playing 1000
air_conditioner 1000
street_music 1000
engine_idling 1000
jackhammer 1000
drilling 1000
siren 929
car_horn 429
gun_shot 374
Name: class, dtype: int64

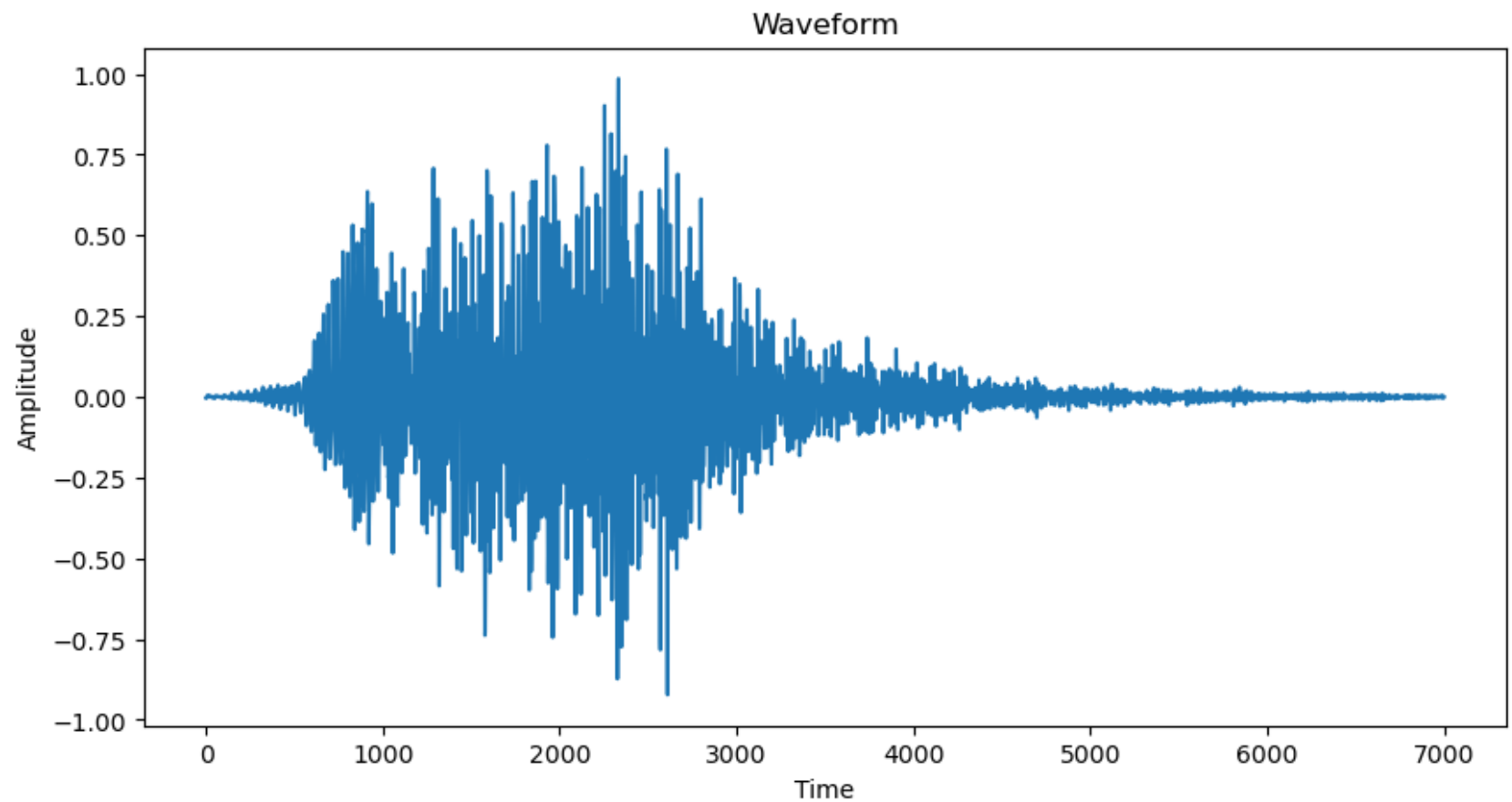
In [17]: 1 len(metadata)

Out[17]: 8732

```
In [18]: 1 metadata.shape
```

```
Out[18]: (8732, 8)
```

```
In [19]: 1 filename2 = 'dog bark2.wav'
2 data, sample_rate = librosa.load(filename2)
3 plt.figure(figsize=(10, 5))
4 plt.plot(data)
5 plt.title('Waveform')
6 plt.xlabel('Time')
7 plt.ylabel('Amplitude')
8 plt.show()
9
10 ipd.Audio(filename2)
```



Out[19]:

0:00 / 0:00

Preprocessing data

```
In [20]: 1 audio_file_path = '100263-2-0-3.wav'  
2 librosa_audio_data, librosa_sample_rate = librosa.load(audio_file_path)
```

```
In [21]: 1 print(librosa_audio_data)
```

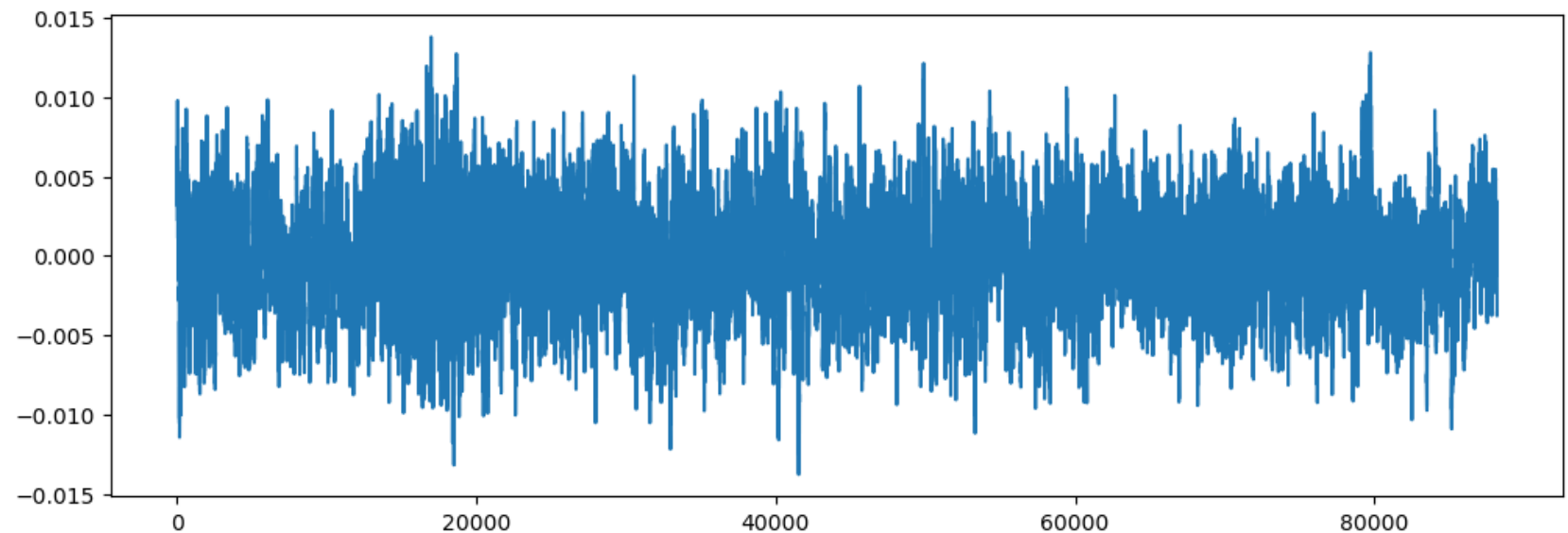
```
[ 0.00331575  0.00467553  0.00361099 ... -0.00376796 -0.00347471  
-0.00357828]
```

Observation

Here Librosa converts the signal to mono, meaning the channel array will be 1:

```
[ 0.00331575 0.00467553 0.00361099 ... -0.00376796 -0.00347471 -0.00357828]
```

```
In [22]: 1 # Let plot the librosa audio data  
2 # original audio with 1 channel  
3 plt.figure(figsize = (12, 4))  
4 plt.plot(librosa_audio_data);
```



```
In [23]: 1 # Let's read with scipy
         2 wave_sample_rate, wave_audio = wav.read(audio_file_path)
```

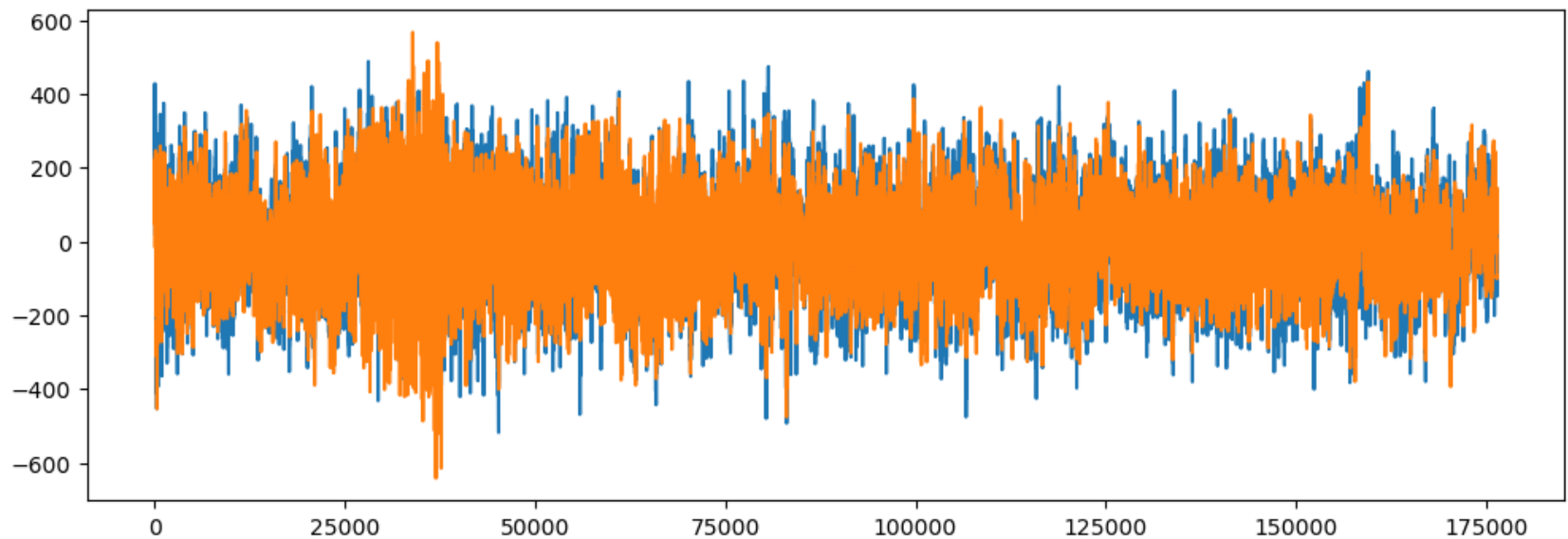
Here we have 2 channels

(array([[194, 100], [179, 113], [160, 124], ..., [-143, -87], [-134, -91], [-110, -98]], dtype=int16))

```
In [24]: 1 wave_audio, wave_sample_rate
```

```
Out[24]: (array([[ 194, 100],
                  [ 179, 113],
                  [ 160, 124],
                  ...,
                  [-143, -87],
                  [-134, -91],
                  [-110, -98]], dtype=int16),
         44100)
```

```
In [25]: 1 # plot Original audio with 2 channels
         2
         3 plt.figure(figsize=(12, 4))
         4 plt.plot(wave_audio);
```



Extracting Features

Here we will be using Mel_Frequency Cepstral Coefficient(MFCC) from the audio samples. The MFCC summarises the frequency distribution across the window size, so it is possible to analyse both the frequency and time characteristics of the sound. These audio representation will allow us to identify features for classification.

<https://chat.openai.com/share/53e7a8f8-90fc-489c-9b4e-e6584c1b02c5> (<https://chat.openai.com/share/53e7a8f8-90fc-489c-9b4e-e6584c1b02c5>).

```
In [26]: 1 mfccs = librosa.feature.mfcc(y = librosa_audio_data, sr = librosa_sample_rate, n_mfcc=40)
```

```
In [27]: 1 mfccs
```

```
Out[27]: array([[ -4.7486273e+02, -4.5088608e+02, -4.4905338e+02, ...,  
                -4.7676157e+02, -4.7334869e+02, -4.9085260e+02],  
               [ 1.1530264e+02,  1.1144249e+02,  1.1125224e+02, ...,  
                1.1112500e+02,  1.1057970e+02,  1.0299150e+02],  
               [-1.8326149e+01, -2.4682453e+01, -3.0259779e+01, ...,  
                -8.2357616e+00, -9.0665359e+00, -4.5019574e+00],  
               ...,  
               [-2.8760438e+00, -3.2479770e+00, -4.8965530e+00, ...,  
                -5.2023709e-01,  3.5672512e+00,  7.4937592e+00],  
               [-4.2968428e-01, -5.8838773e-01, -8.1724101e-01, ...,  
                1.8340731e-01,  7.6732409e-01,  2.7120016e+00],  
               [-1.1780634e+00,  6.9809389e-01,  6.3521605e+00, ...,  
                -2.6221924e+00, -4.7912717e+00, -3.1826315e+00]], dtype=float32)
```

```
In [28]: 1 print(mfccs.shape)
```

```
(40, 173)
```

Extracting MFCC's for every audio file

```
In [29]: 1 # pip install resampy
        2
```

```
In [30]: 1 import os
        2 import numpy as np
        3
```

```
In [31]: 1 audio_dataset_path = 'UrbanSound8K/audio/'
        2 metadata = pd.read_csv('UrbanSound8K/metadata/UrbanSound8K.csv')
        3 metadata.head()
```

Out[31]:

	slice_file_name	fsID	start	end	salience	fold	classID	class
0	100032-3-0-0.wav	100032	0.0	0.317551	1	5	3	dog_bark
1	100263-2-0-117.wav	100263	58.5	62.500000	1	5	2	children_playing
2	100263-2-0-121.wav	100263	60.5	64.500000	1	5	2	children_playing
3	100263-2-0-126.wav	100263	63.0	67.000000	1	5	2	children_playing
4	100263-2-0-137.wav	100263	68.5	72.500000	1	5	2	children_playing

```
In [32]: 1 import os
        2 import librosa
        3 import numpy as np
        4 from tqdm import tqdm
        5
        6 # Creating a function to extract scaled MFCC features
        7 def features_extractor(file_name):
        8     audio, sample_rate = librosa.load(file_name, res_type='kaiser_fast')
        9     mfccs_features = librosa.feature.mfcc(y=audio, sr=sample_rate, n_mfcc=40)
       10     mfccs_scaled_features = np.mean(mfccs_features.T, axis=0)
       11     return mfccs_scaled_features
```


In [82]: 1 `len(extracted_features)`

Out[82]: 8732

In [83]: 1 `extracted_features_df = pd.DataFrame(extracted_features, columns=['feature', 'class'])`
 2 `extracted_features_df.head()`

Out[83]:

	feature	class
0	[-217.35526, 70.22338, -130.38527, -53.282898, ...]	dog_bark
1	[-424.09818, 109.34077, -52.919525, 60.86475, ...]	children_playing
2	[-458.79114, 121.38419, -46.52066, 52.00812, ...]	children_playing
3	[-413.89984, 101.66371, -35.42945, 53.036354, ...]	children_playing
4	[-446.60352, 113.68541, -52.402214, 60.302044, ...]	children_playing

In [84]: 1 *# Split the dataset into independent and dependednt dataset*
 2 `x = np.array(extracted_features_df['feature'].tolist())`
 3 `y = np.array(extracted_features_df['class'].tolist())`

In [85]: 1 `x`

Out[85]: array([[-2.1735526e+02, 7.0223381e+01, -1.3038527e+02, ...,
 -1.6930529e+00, -6.1698329e-01, 3.8600507e-01],
 [-4.2409818e+02, 1.0934077e+02, -5.2919525e+01, ...,
 5.3489321e-01, -5.4468715e-01, 4.4632098e-01],
 [-4.5879114e+02, 1.2138419e+02, -4.6520660e+01, ...,
 2.0768483e+00, 1.6962965e+00, -9.6140963e-01],
 ...,
 [-3.0388824e+02, 1.1135945e+02, -4.5941566e+01, ...,
 -3.0292377e+00, 2.7170296e+00, 7.6197419e+00],
 [-3.4411008e+02, 1.2545021e+02, -5.4903442e+01, ...,
 -7.9082427e+00, -1.6414584e+00, 5.6668439e+00],
 [-3.1560281e+02, 9.4854805e+01, -3.7222340e+01, ...,
 6.1386460e-01, -1.1449189e+01, -6.0105853e+00]], dtype=float32)

In [86]: 1 y

Out[86]: array(['dog_bark', 'children_playing', 'children_playing', ...,
 'car_horn', 'car_horn', 'car_horn'], dtype='<U16')

In [87]: 1 x.shape, y.shape

Out[87]: ((8732, 40), (8732,))

```
In [89]: 1 # Label encoding (converting it ti numbers)
2 # y = np.array(pd.get_dummies(y))
3 from tensorflow.keras.utils import to_categorical
4 from sklearn.preprocessing import LabelEncoder
5 labelencoder = LabelEncoder()
6 y = to_categorical(labelencoder.fit_transform(y))
```

In [90]: 1 y

Out[90]: array([[0., 0., 0., ..., 0., 0., 0.],
 [0., 0., 1., ..., 0., 0., 0.],
 [0., 0., 1., ..., 0., 0., 0.],
 ...,
 [0., 1., 0., ..., 0., 0., 0.],
 [0., 1., 0., ..., 0., 0., 0.],
 [0., 1., 0., ..., 0., 0., 0.]], dtype=float32)

In [91]: 1 y.shape

Out[91]: (8732, 10)

```
In [92]: 1 from sklearn.model_selection import train_test_split
2 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=0)
```

In [93]: 1 x_train

```
Out[93]: array([[ -1.3110471e+02,  1.1250591e+02, -2.2574696e+01, ...,
        3.2466526e+00, -1.3690237e+00,  2.7557542e+00],
       [ -1.3670342e+01,  9.1085083e+01, -7.7927337e+00, ...,
        -3.2530508e+00, -5.2774529e+00, -1.5569715e+00],
       [ -4.9871544e+01,  2.6535299e-01, -2.0500937e+01, ...,
        2.8545945e+00, -1.6092046e+00,  3.5248058e+00],
       ...,
       [ -4.2701236e+02,  9.2623047e+01,  3.1293974e+00, ...,
        7.4264139e-01,  7.3349088e-01,  7.1100914e-01],
       [ -1.4575461e+02,  1.3626578e+02, -3.3515518e+01, ...,
        1.4681193e+00, -2.0091701e+00, -8.8218188e-01],
       [ -4.2103134e+02,  2.1065454e+02,  3.4906609e+00, ...,
        -5.3888674e+00, -3.3713605e+00, -1.5665115e+00]], dtype=float32)
```

In [94]: 1 x_train.shape, y_train.shape, x_test.shape, y_test.shape

```
Out[94]: ((6985, 40), (6985, 10), (1747, 40), (1747, 10))
```

Model creation

```
In [95]: 1 import tensorflow as tf
        2 import keras
        3
        4 print("TensorFlow version:", tf.__version__)
        5 print("Keras version:", keras.__version__)
```

TensorFlow version: 2.12.0

Keras version: 2.12.0

```
In [96]: 1 import tensorflow as tf
        2 from tensorflow.keras.models import Sequential
        3 from tensorflow.keras.layers import Dense, Dropout, Activation, Flatten
        4 from tensorflow.keras.optimizers import Adam
        5 from sklearn import metrics
        6
```

```
In [97]: 1 # number of classes  
2 num_labels = y.shape[1]
```

```
In [98]: 1 Dense(units=32, activation='relu')  
2
```

Out[98]: <keras.layers.core.dense.Dense at 0x1ffb539dc40>

```
In [99]: 1 model = Sequential()  
2  
3 # first layer  
4 model.add(Dense(100, input_shape=(40,)))  
5 model.add(Activation('relu'))  
6 model.add(Dropout(0.5))  
7  
8 # second layer  
9 model.add(Dense(200))  
10 model.add(Activation('relu'))  
11 model.add(Dropout(0.5))  
12  
13 # third layer  
14 model.add(Dense(100))  
15 model.add(Activation('relu'))  
16 model.add(Dropout(0.5))  
17  
18 # final layer  
19 model.add(Dense(num_labels))  
20 model.add(Activation('softmax'))
```

In [100]: 1 model.summary()

Model: "sequential_1"

Layer (type)	Output Shape	Param #
dense_6 (Dense)	(None, 100)	4100
activation_4 (Activation)	(None, 100)	0
dropout_3 (Dropout)	(None, 100)	0
dense_7 (Dense)	(None, 200)	20200
activation_5 (Activation)	(None, 200)	0
dropout_4 (Dropout)	(None, 200)	0
dense_8 (Dense)	(None, 100)	20100
activation_6 (Activation)	(None, 100)	0
dropout_5 (Dropout)	(None, 100)	0
dense_9 (Dense)	(None, 10)	1010
activation_7 (Activation)	(None, 10)	0
Total params: 45,410		
Trainable params: 45,410		
Non-trainable params: 0		

In [101]: 1 model.compile(optimizer='adam', loss='categorical_crossentropy', metrics=['accuracy'])


```
In [102]: 1 # Training our model
2 from tensorflow.keras.callbacks import ModelCheckpoint
3 from datetime import datetime
4
5 num_epochs = 100
6 num_batch_size = 32
7
8 checkpointer = ModelCheckpoint(filepath='model/audio_classification.hdf5', verbose=1, save_best_only=
9 start = datetime.now()
10
11 model.fit(x_train, y_train, batch_size=num_batch_size, epochs=num_epochs, validation_data=(x_test, y_
12
13 duration = datetime.now() - start
14 print('Training completed in time:', duration)
15
```

210/219 [=====>.] - ETA: 0s - loss: 0.9477 - accuracy: 0.6851
Epoch 97: val_loss improved from 0.74097 to 0.73980, saving model to model\audio_classification.hdf5
219/219 [=====] - 2s 8ms/step - loss: 0.9511 - accuracy: 0.6842 - val_loss: 0.7398 - val_accuracy: 0.7705
Epoch 98/100
212/219 [=====>.] - ETA: 0s - loss: 0.9590 - accuracy: 0.6891
Epoch 98: val_loss improved from 0.73980 to 0.73447, saving model to model\audio_classification.hdf5
219/219 [=====] - 2s 9ms/step - loss: 0.9584 - accuracy: 0.6885 - val_loss: 0.7345 - val_accuracy: 0.7768
Epoch 99/100
218/219 [=====>.] - ETA: 0s - loss: 0.9358 - accuracy: 0.6892
Epoch 99: val_loss did not improve from 0.73447
219/219 [=====] - 1s 6ms/step - loss: 0.9358 - accuracy: 0.6892 - val_loss: 0.7430 - val_accuracy: 0.7687
Epoch 100/100
212/219 [=====>.] - ETA: 0s - loss: 0.9290 - accuracy: 0.6862
Epoch 100: val_loss did not improve from 0.73447
219/219 [=====] - 1s 6ms/step - loss: 0.9281 - accuracy: 0.6875 - val_loss: 0.7567 - val_accuracy: 0.7710
Training completed in time: 0:02:47.064258

```
In [103]: 1 test_accuracy = model.evaluate(x_test, y_test, verbose=0)
2 test_accuracy[1]
```

Out[103]: 0.7710360884666443

```
In [104]: 1 # Extract features from the audio file
2 filename = 'dog_bark.wav'
3 prediction_feature = features_extractor(filename)
4 prediction_feature = prediction_feature.reshape(1, -1)
5
6 # Make predictions
7 predictions = model.predict(prediction_feature)
8 predicted_classes = np.argmax(predictions, axis=-1)
9
10 print(predicted_classes)
```

```
1/1 [=====] - 0s 110ms/step
[3]
```

```
In [113]: 1 metadata
```

Out[113]:

	slice_file_name	fsID	start	end	salience	fold	classID	class
0	100032-3-0-0.wav	100032	0.000000	0.317551	1	5	3	dog_bark
1	100263-2-0-117.wav	100263	58.500000	62.500000	1	5	2	children_playing
2	100263-2-0-121.wav	100263	60.500000	64.500000	1	5	2	children_playing
3	100263-2-0-126.wav	100263	63.000000	67.000000	1	5	2	children_playing
4	100263-2-0-137.wav	100263	68.500000	72.500000	1	5	2	children_playing
...
8727	99812-1-2-0.wav	99812	159.522205	163.522205	2	7	1	car_horn
8728	99812-1-3-0.wav	99812	181.142431	183.284976	2	7	1	car_horn
8729	99812-1-4-0.wav	99812	242.691902	246.197885	2	7	1	car_horn
8730	99812-1-5-0.wav	99812	253.209850	255.741948	2	7	1	car_horn
8731	99812-1-6-0.wav	99812	332.289233	334.821332	2	7	1	car_horn

8732 rows × 8 columns

```
In [106]: 1 # Filter metadata based on predicted class ID
2 predicted_class_id = predicted_classes[0]
3 filtered_metadata = metadata[metadata['classID'] == predicted_class_id]
4 # Load metadata
5
6 metadata = pd.read_csv('UrbanSound8K/metadata/UrbanSound8K.csv')
7
8 if len(filtered_metadata) == 0:
9     print("Predicted class ID not found in the metadata.")
10 else:
11     # Extract relevant columns from filtered metadata
12     filtered_metadata = filtered_metadata[['classID', 'class']]
13
14     # Display predicted class and class ID
15     prediction = filtered_metadata.iloc[0]
16     print(f"Predicted class: {prediction['class']}")
17     print(f"Predicted class ID: {prediction['classID']}")
18
```

Predicted class: dog_bark

Predicted class ID: 3

In [107]:

```
1 import pandas as pd
2 from tabulate import tabulate
3
4 # Filter metadata based on predicted class ID
5 predicted_class_id = predicted_classes[0]
6 filtered_metadata = metadata[metadata['classID'] == predicted_class_id]
7
8 # Load metadata
9 metadata = pd.read_csv('UrbanSound8K/metadata/UrbanSound8K.csv')
10
11 if len(filtered_metadata) == 0:
12     print("Predicted class ID not found in the metadata.")
13 else:
14     # Extract relevant columns from filtered metadata
15     filtered_metadata = filtered_metadata[['classID', 'class']]
16
17     # Display first 10 rows of the filtered metadata in a table
18     table_data = filtered_metadata.head(10).values.tolist()
19     headers = filtered_metadata.columns.tolist()
20     print(tabulate(table_data, headers, tablefmt='grid'))
21
```

classID	class
3	dog_bark
3	dog_bark
3	dog_bark
3	dog_bark
3	dog_bark
3	dog_bark
3	dog_bark
3	dog_bark
3	dog_bark
3	dog_bark
3	dog_bark
3	dog_bark

Testing some audio data

- Preprocess the new audio data
- predict the class
- inverse transform your predicted label

In [114]:

```
1  # Specify the file path
2  filename = '21684-9-0-5.wav'
3
4  # Check if the file exists
5  try:
6      open(filename)
7  except FileNotFoundError:
8      print(f"File '{filename}' not found.")
9      exit()
10
11 # Load the audio file
12 audio, sample_rate = librosa.load(filename, res_type='kaiser_fast')
13
14 # Extract MFCC features
15 mfccs_features = librosa.feature.mfcc(y=audio, sr=sample_rate, n_mfcc=40)
16
17 # Scale the features
18 mfccs_scaled_features = np.mean(mfccs_features.T, axis=0)
19
20 # Print the scaled features
21 print(mfccs_scaled_features)
22
23 # Reshape the features
24 mfccs_scaled_features = mfccs_scaled_features.reshape(1, -1)
25
26 # Print the reshaped features
27 print(mfccs_scaled_features)
28 print(mfccs_scaled_features.shape)
29
30 # Make predictions using the model
31 predicted_probabilities = model.predict(mfccs_scaled_features)
32
33 # Get the predicted label with the highest probability
34 predicted_label = np.argmax(predicted_probabilities)
35
36 # Print the predicted label
37 print(predicted_label)
38
39 # Decode the predicted label using the Labelencoder
40 prediction_class = labelencoder.inverse_transform([predicted_label])
41
42 # Print the predicted class
```

```
43 print(prediction_class)
```

```
[ -178.82379      136.7633      -13.068221      48.14213       6.6529574
   36.10647      -12.541818      -3.522911      -12.262587      10.594036
   -4.1306596      14.425096      -2.3365233      16.041061       0.95959246
   13.80531       -1.9968798       5.126442      -2.5143142       4.8222094
   -6.7894983       2.2693586      -0.7009771       2.7691789      -0.5069647
   -0.89399356      -1.6918974       3.7261329       2.4733682      -1.0054693
   -3.8394966       0.2816141       6.8972626       2.6243734      -2.8658035
   -3.4202144      -1.4473431      -9.400884      -1.6816708       2.3436675 ]
[[ -178.82379      136.7633      -13.068221      48.14213       6.6529574
   36.10647      -12.541818      -3.522911      -12.262587      10.594036
   -4.1306596      14.425096      -2.3365233      16.041061       0.95959246
   13.80531       -1.9968798       5.126442      -2.5143142       4.8222094
   -6.7894983       2.2693586      -0.7009771       2.7691789      -0.5069647
   -0.89399356      -1.6918974       3.7261329       2.4733682      -1.0054693
   -3.8394966       0.2816141       6.8972626       2.6243734      -2.8658035
   -3.4202144      -1.4473431      -9.400884      -1.6816708       2.3436675 ]]
(1, 40)
1/1 [=====] - 0s 42ms/step
9
['street_music']
```

In []:

1